

# DAT105 Assignment 3 - Group 2

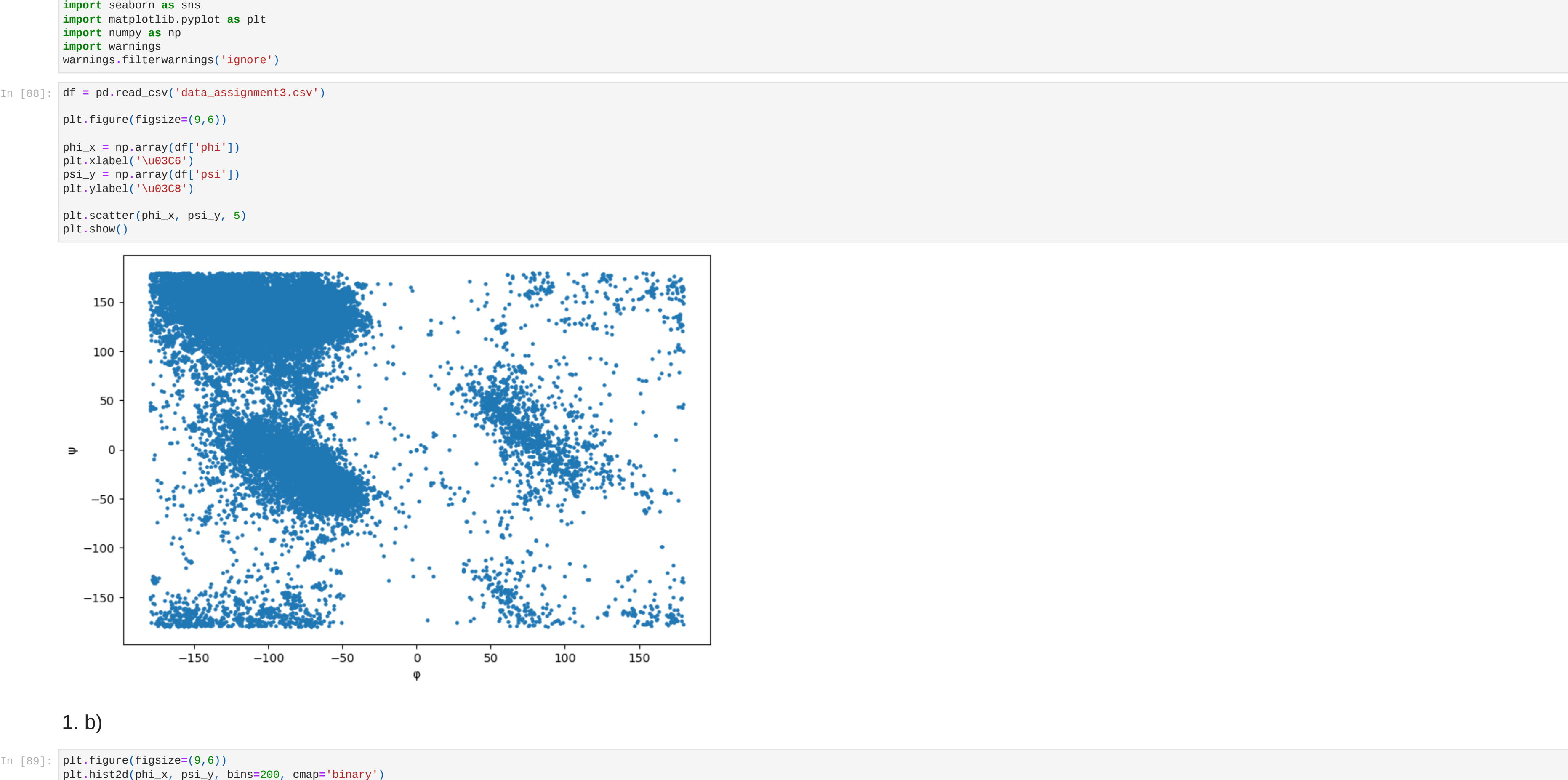
Francisco Boudagh - (15 hours)

Jakob Engström - (15 hours)

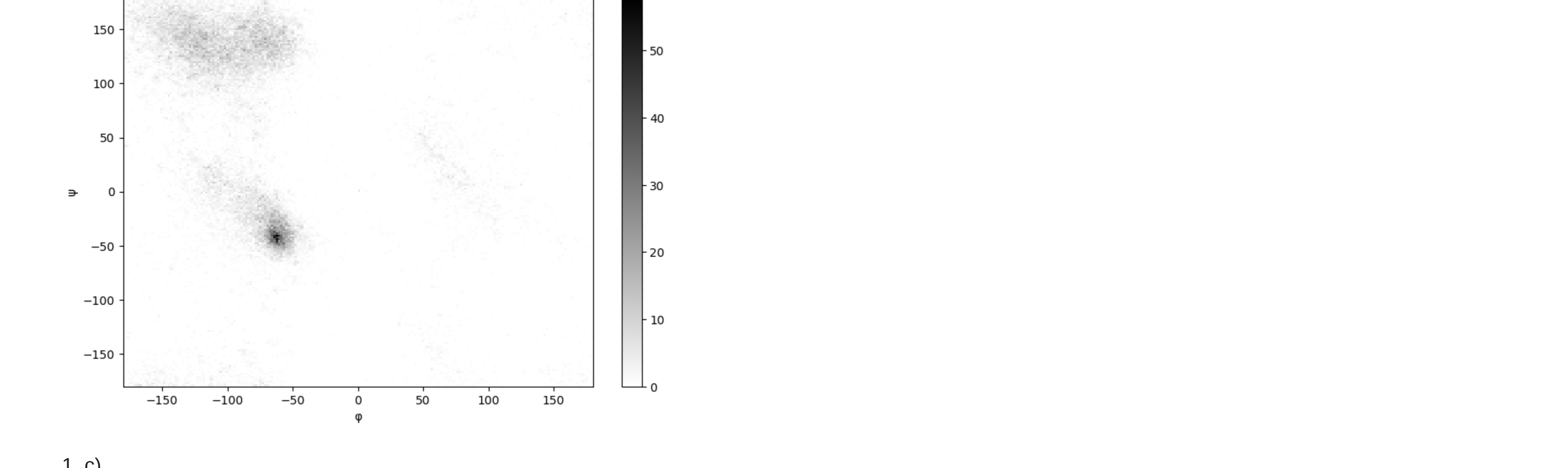
April 19, 2023

## Problem 1

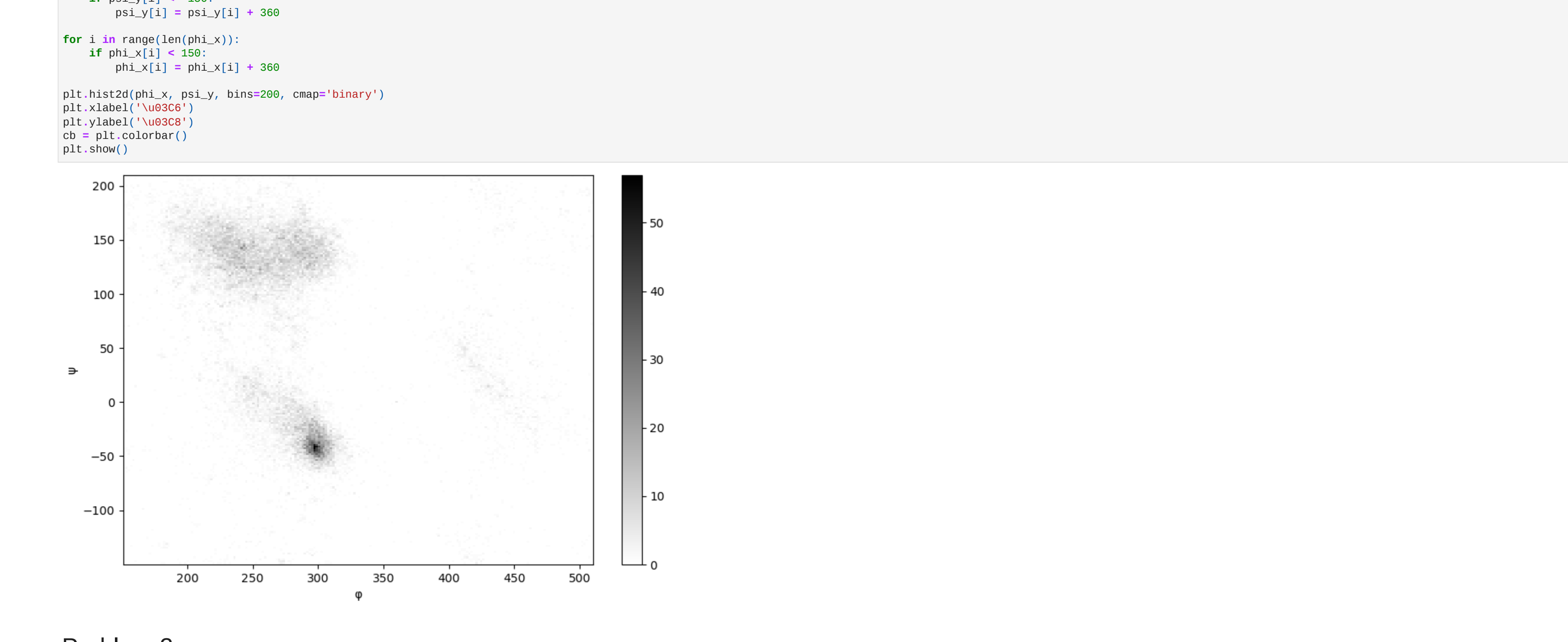
### 1. a)



### 1. b)

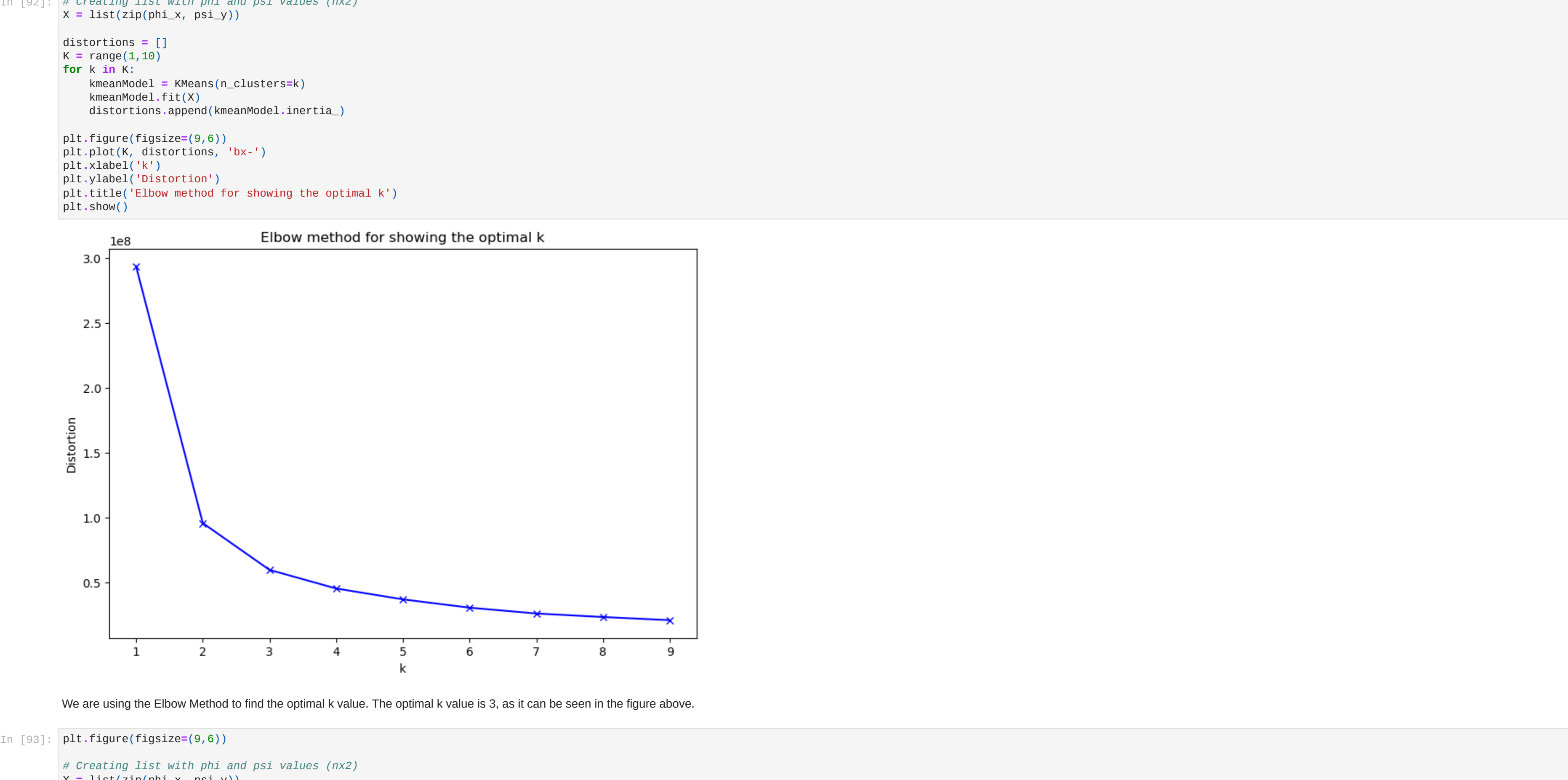


### 1. c)

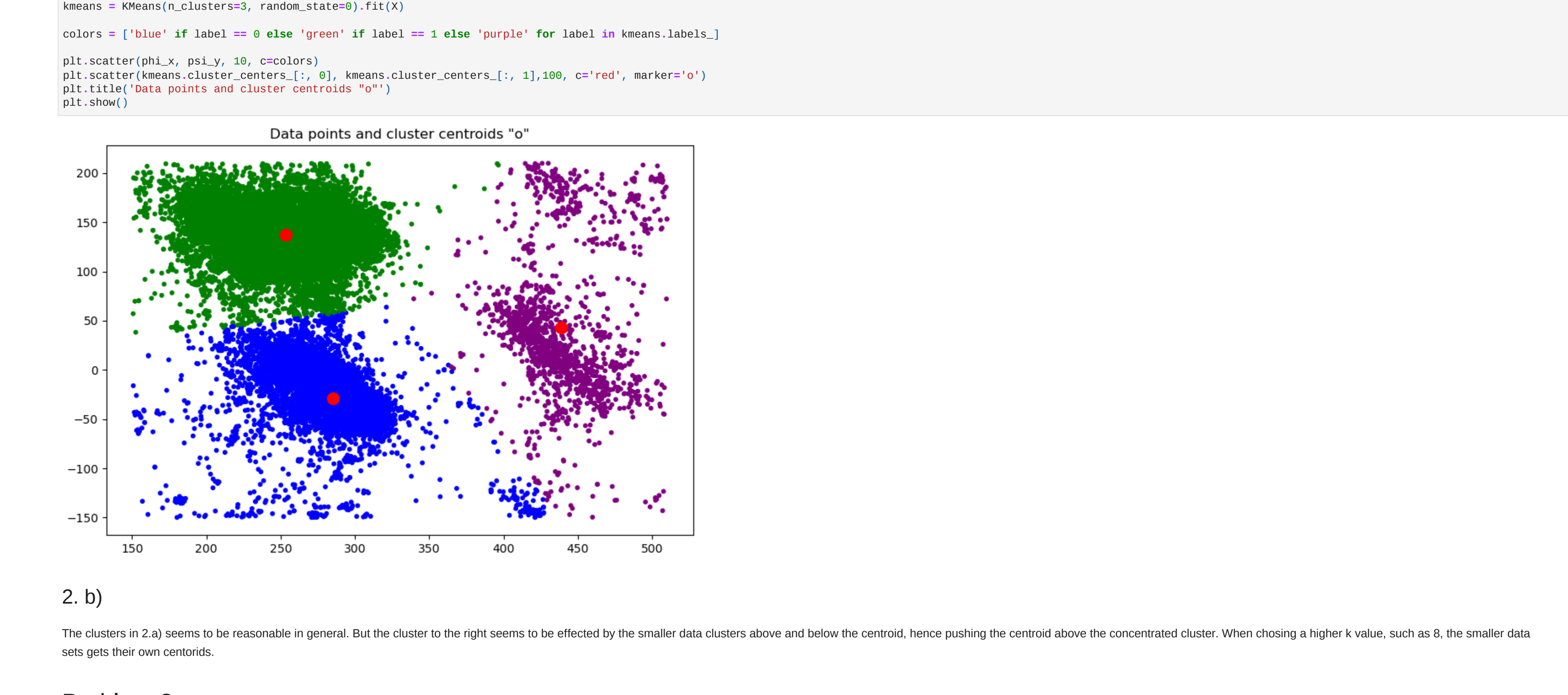


## Problem 2

### 2. a)



We are using the Elbow Method to find the optimal k value. The optimal k value is 3, as it can be seen in the figure above.

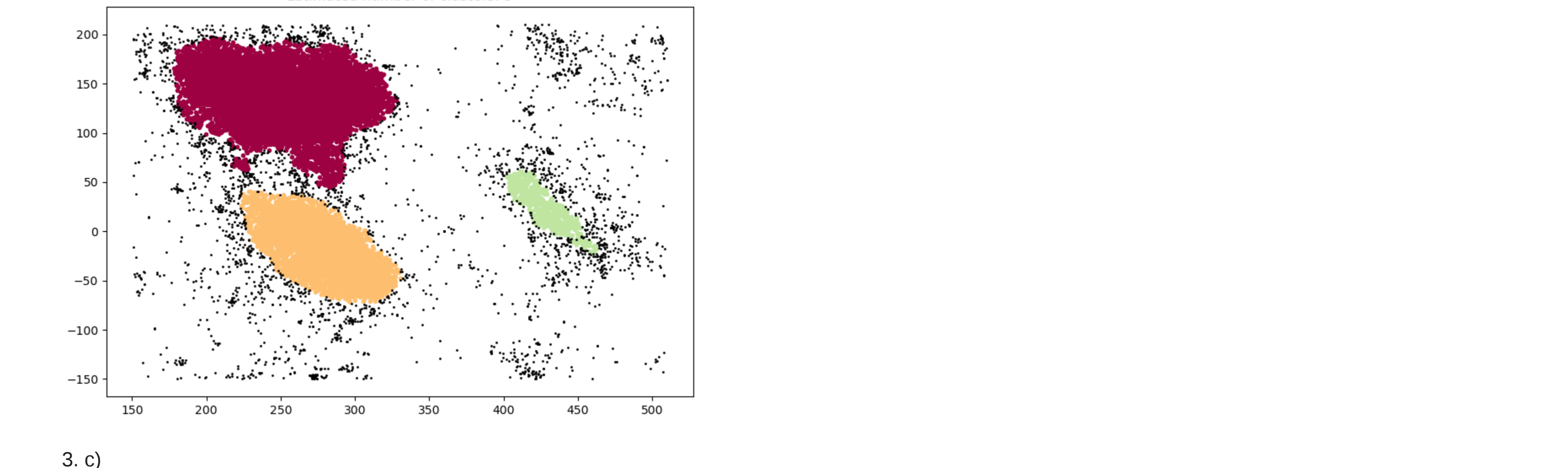
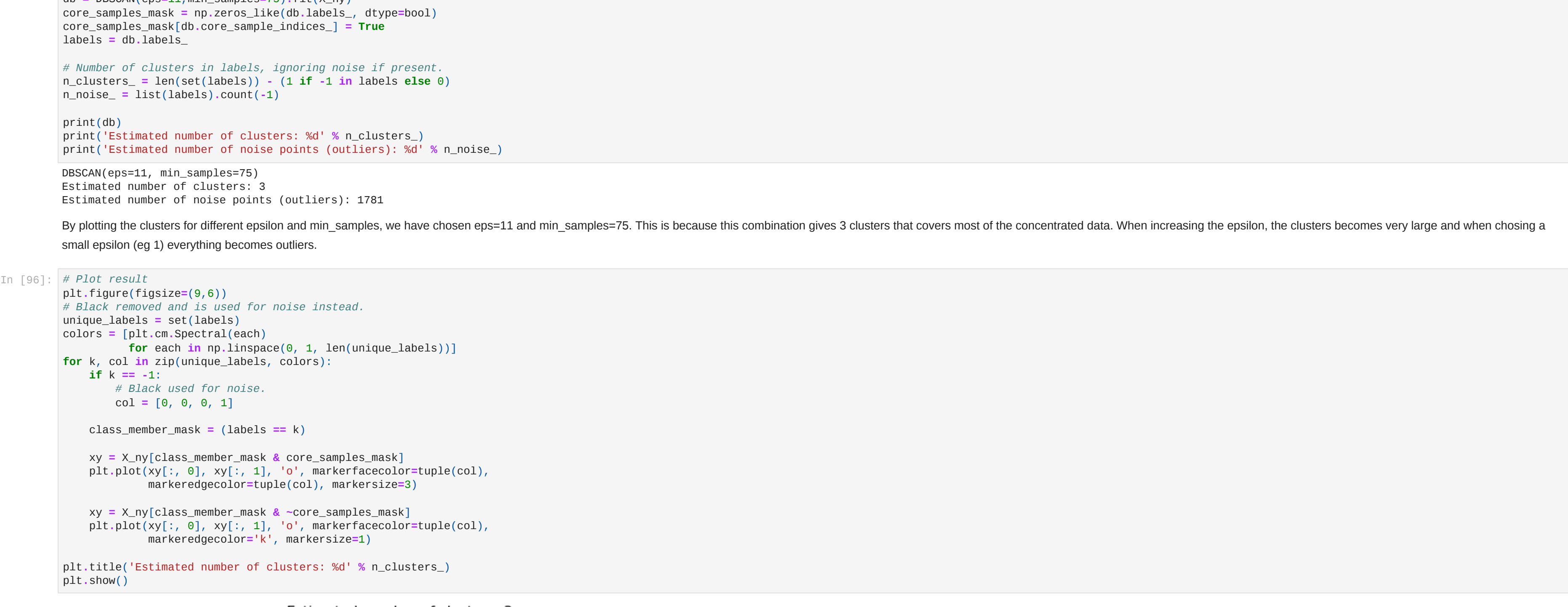


### 2. b)

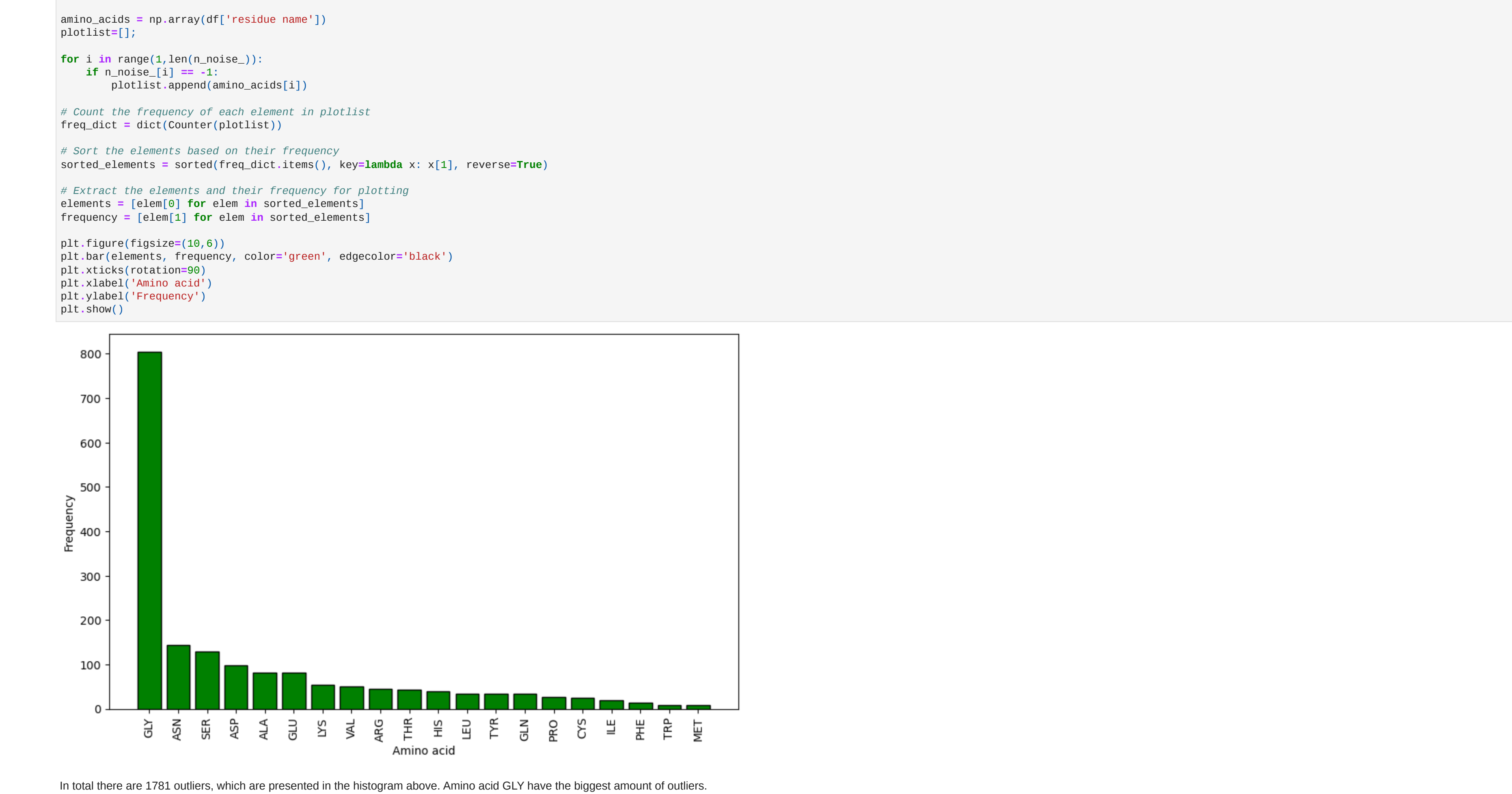
The clusters in 2.a) seems to be reasonable in general. But the cluster to the right seems to be effected by the smaller data clusters above and below the centroid, hence pushing the centroid above the concentrated cluster. When choosing a higher k value, such as 8, the smaller data sets gets their own centroids.

## Problem 3

### 3. a) and b)



### 3. c)



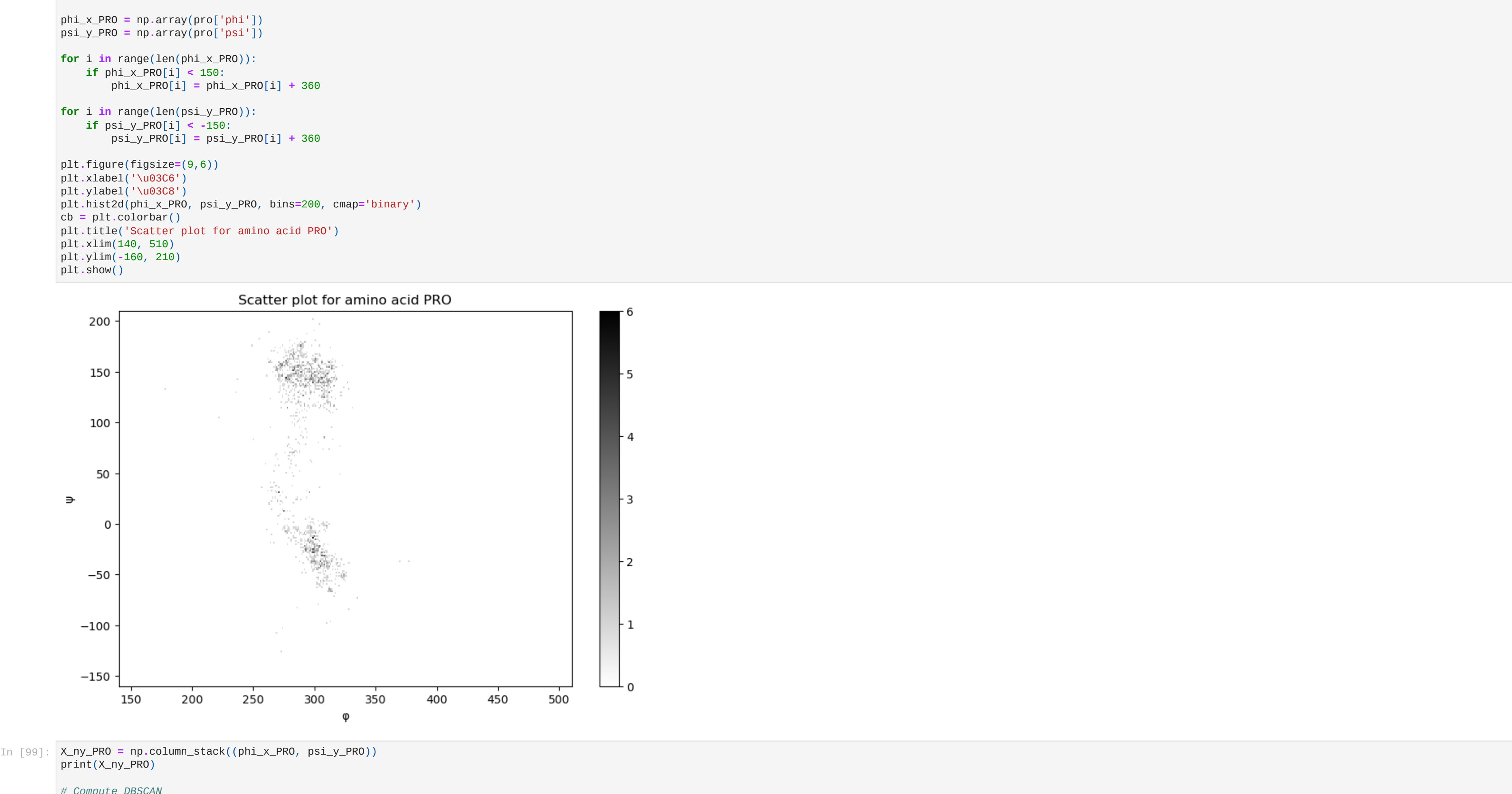
In total there are 1781 outliers, which are presented in the histogram above. Amino acid GLY have the biggest amount of outliers.

### 3. d)

Comparison between K-means and DBSCAN  
When using K-means, we include outliers in the clusters, which can make K-means sensitive to outliers and affect the clusters. Additionally, K-means works well for circular datasets, but in this particular dataset, especially the cluster to the right is not very spherical, so K-means is probably not the optimal choice. On the other hand, when using DBSCAN, we obtain a reasonable cluster for the dataset because this method does not include outliers and can create clusters for almost any kind of dataset, regardless of its shape (i.e., it does not need to be spherical).

## Problem 4

### Scatter of psi-phi for amino acid 'PRO'



We have transformed the data for amino acid PRO and set x and y axis limits so we can see where the PRO data were placed in the DBSCAN plot in 3.b).

From the DBSCAN above, we can see that there are mainly two variants of the amino acid PRO. Both have mainly the same  $\phi$  angle (around 300°) but two different  $\psi$  angles (at -40° and 150°). The plot above is better when looking at a specific amino acid (in this case PRO), compared to the DBSCAN plot in 3.b) where we are plotting and using DBSCAN to cluster all the amino acids. In the plot from 3.b) we can not identify different amino acids just by looking at the plot, since the plot in 3.b) is just showing amino acids with similar  $\phi$ - $\psi$  angles combination.